

# Large deviation theory to model systems under an external feedback

Alessio Gagliardi<sup>1</sup>, Alessandro Pecchia<sup>2</sup>, Aldo Di Carlo<sup>3</sup>

(1) *Technische Universitaet Muenchen,  
Arcisstrasse 21, 80333, Munich (Germany)  
alessio.gagliardi@tum.de.*

(2) *CNR, Via Salaria Km 29,  
600, 0017 Monte Rotondo (Italy)*

(3) *University of Rome "Tor Vergata",  
Via del Politecnico 1, 00133, Rome (Italy)*

(Dated: April 18, 2016)

In this paper we address the problem of systems under an external feedback. This is performed using a large deviation approach and rate distortion from information theory. In particular we define a lower boundary for the maximum entropy reduction that can be obtained using a feedback apparatus with a well defined accuracy in terms of measurement of the state of the system. The large deviation approach allows also to define a new set of potentials, including information, which similarly to more conventional thermodynamic potentials can define the state with optimal use of the information given the accuracy of the feedback apparatus.

**PACS numbers:**

## I. INTRODUCTION

The idea of a thermodynamic theory of systems under an external feedback dates back to the origin of statistical physics when Maxwell made a gedanken experiment about the work that could be extracted by a system controlled by an external apparatus. In principle the external feedback can probe the state of the system and manipulates it in order to extract work. In his first formulation the feedback was treated like an oracle that can make measures and manipulate the system without any cost in terms of work and entropy production, for this reason it was called a "Demon" to stress its abstract nature. However the idea was reconsidered along all the past century<sup>1,2</sup> with different approaches and conclusions. One of the central work on the topic was the paper by Szilard and his famous Szilard engine<sup>3</sup>, a still very idealized machine that anyway shows already some more realistic characteristics in terms of its components and the nature of the feedback controller. Starting from Szilard engine, but generalizing the concept, several authors<sup>4-10</sup> and in particular Sagawa and Ueda, have developed a new branch of thermodynamics of systems under an external feedback<sup>11-13</sup>, i.e. information thermodynamics.

In particular in<sup>14</sup> it was shown, using an ingenious extension of fluctuation theorems including the information gathered by the feedback, that the presence of the external controller can increase the average work that can be extracted from a system when it is driven from one equilibrium state to another according to

$$\langle W \rangle \leq \Delta F + k_B T I, \quad (1)$$

with  $\langle W \rangle$  the average work,  $\Delta F$  free energy variation between the final and initial equilibrium states,  $k_B$  is the Boltzmann constant,  $T$  the bath temperature and  $I$  the mutual information. For a cyclic transformation, for which  $\Delta F = 0$  the maximum work reduces to  $k_B T I$ .

The mutual information, for a classical system, is described by the mutual information functional of information theory. There is an understandable debate about the validity of using information theoretical quantities and their interpretations in thermodynamics<sup>15</sup>. However, for many cases of interest<sup>16</sup> Shannon entropy represents a good choice for the entropy functional, especially in equilibrium cases. In particular Shannon entropy gives the right entropy functional but only physical and experimental evidences can determine which is the correct probability density function (PDF) of the problem under investigation, as information theory has nothing to say about that<sup>17</sup>.

An interesting alternative approach to non equilibrium thermodynamics is arising within the framework of a different branch of information theory, e.g. large deviation theory (LDT). LDT is a mathematical framework used to establish the probability of fluctuations of statistical quantities from their *typical values*. Typicality in information theory leads to the definition of thermodynamic averages and state variables<sup>18,19</sup>.

In particular, LDT demonstrates that for a broad class of PDFs, the probability of fluctuations from the average value drops exponentially with the fluctuation magnitude. The exponent is proportional to the number of degrees of freedom of the system, explaining why, in the thermodynamic limit, fluctuations of state variables become negligible. However in the emerging field of stochastic thermodynamics which also copes with small systems in non equilibrium conditions, fluctuations can be an important aspect of their behavior<sup>20-24</sup>.

Many systems of interest have the peculiarity of being the interconnection between a physical system probed and controlled by a feedback apparatus, for example biological processes in cells and other living organisms belong to this class. Recently, several studies have investigated the effect of including information terms in the analysis of biological processes within the framework of

information thermodynamics with great success, see for example<sup>25–32</sup>. Thus the investigation of thermodynamics under feedback is rapidly rising interest in many fields beyond statistical physics.

In this paper we give a new insight to the problem of a system under an external feedback using a special case of LDT and the concept of typicality. The link between systems under feedback and LDT is obtained using rate distortion theory (RDT), a fundamental part of information theory. The chain of relationships between LDT and RDT that we outline, is particularly pleasing as it makes a direct connection between the information the feedback controller apparatus gathers about a system and the associated entropy reduction, also leading to an explicit construction of a thermodynamic potential for a system under feedback. The possibility to construct potentials including the effect of information opens the perspective of a complete new analysis of those systems more similar to conventional equilibrium thermodynamic formalism. For a similar approach see also<sup>33</sup>.

The paper is organized as follows: in the first part a brief summary of the main information theoretical concepts, Shannon entropy, conditional entropy and mutual information is given. Then the concept of typicality is introduced. In the next section we present the large deviation theory and the code large distortion problem. The final part is devoted to introducing rate distortion theory and the final link to thermodynamics under an external feedback.

## II. A BRIEF EXCURSUS IN INFORMATION THEORY

The central quantity of information theory is the Shannon entropy defined as

$$S(P) = -k_B \sum_i p_i \ln p_i, \quad (2)$$

where  $p_i$  represents the probability of the  $i^{th}$  event and  $P$  the PDF. Shannon entropy is usually expressed in bits and adimensional, we have chosen here the convention to include the Boltzmann constant and express the entropy using natural logarithm. However, this is totally immaterial for the present discussion. Shannon entropy plays a central role in thermodynamics, even for a large class of systems under non equilibrium conditions<sup>16</sup>.

More generally if we have a PDF  $p(\xi)$  which depends on a set of degrees of freedom,  $\xi$ , we can define the proper discrete (or continuous) Shannon entropy. The degrees of freedom depend on the problem at hand, they could be the set of positions and momenta of a collection of particle in gas phase for example. From the Shannon entropy it is possible to promptly derive other four connected quantities. The first is the joint (Shannon) entropy which is just the entropy of two random vectors  $\xi$  and  $\pi$ :

$$S(\Xi, \Pi) = -k_B \sum p(\xi, \pi) \ln p(\xi, \pi), \quad (3)$$

with  $p(\xi, \pi)$  the joint probability. If the two vectors are independent the joint entropy reduces to the sum of the individual entropies, in the other case we have:

$$S(\Xi, \Pi) = S(\Xi) + S(\Pi) - k_B I(\Xi, \Pi), \quad (4)$$

with  $I(\Xi, \Pi)$  the mutual information functional defined as:

$$I(\Xi, \Pi) = \sum p(\xi, \pi) \ln \left[ \frac{p(\xi, \pi)}{p(\xi)p(\pi)} \right]. \quad (5)$$

The mutual information is a special case of the Kullback-Leibler divergence (KLd) defined as:

$$D(P||Q) = \sum p(\xi) \ln \left( \frac{p(\xi)}{q(\xi)} \right), \quad (6)$$

where  $P$  and  $Q$  are two PDFs. The KLd is a pseudo distance between PDFs and has an important role in LDT due to the Chernoff bound<sup>19</sup>, in stochastic thermodynamics<sup>34</sup>, within fluctuation theorems<sup>35</sup> and in entropy production within the Boltzmann equation<sup>36</sup>. In particular, we observe that the mutual information has the form of a KLd between the joint PDF,  $p(\xi, \pi)$ , and the independent PDF,  $p(\xi, \pi) = p(\xi)p(\pi) \equiv q(\xi, \pi)$ .

Finally, mutual information and entropy of two random vectors are connected by the conditional entropy:

$$S(\Xi|\Pi) = S(\Xi) - k_B I(\Xi, \Pi) = \sum p(\xi, \pi) \ln p(\xi|\pi). \quad (7)$$

Conditional entropy represents the entropy (uncertainty) left in  $\Xi$  after the conditioning to  $\Pi$ .

These five are the most relevant functionals in information theory as practically every important theorem is related to one or another in some form.

## III. TYPICAL SET IN THE PHASE SPACE

The Shannon entropy has a nice geometrical interpretation in the typical set theorem consequence of the asymptotic equipartition principle<sup>37</sup>. The theorem states that given a PDF  $p(\xi)$ , where  $\xi$  is the vector of degrees of freedom of the problem, with entropy  $S(\Xi)$  then the "typical set" within the phase space has a volume equal to:

$$\Omega_{typ} \sim e^{\frac{S(\Xi)}{k_B}}, \quad (8)$$

where the meaning of "typical" means that the probability for the system to be found in a microstate within the typical set converges to 1 in the thermodynamic limit, namely,

$$p(\xi \in \Omega_{typ}) \rightarrow 1. \quad (9)$$

In other words the concept of typicality states that only a portion of the entire phase space is really relevant to

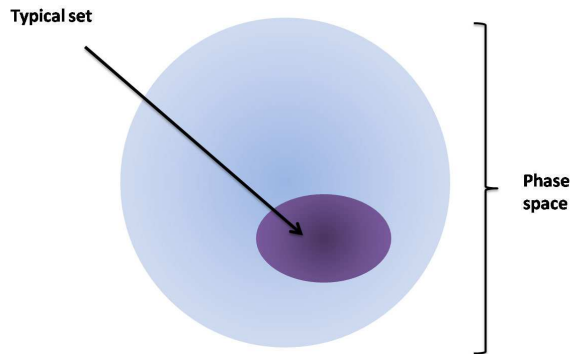


FIG. 1. (color online) The phase space and the typical subset. Usually, except for the uniform distribution, the typical set is indeed a proper subset of the entire class of possible microstates in the phase space. Its volume grows exponentially with the entropy of the problem.

compute ensemble averages of thermodynamic quantities, considering that the volume  $\Omega_{typ}$  fundamentally collects all the probability (see Fig. 1).

In particular for any quantity,  $A(\xi)$ , defined over the phase space, the ensemble average is defined as,

$$\langle A \rangle = \int \tilde{A} q(\tilde{A}) d\tilde{A}, \quad (10)$$

with

$$q(\tilde{A}) = \int p(\xi) \delta(\tilde{A} - A(\xi)) d\xi. \quad (11)$$

If it happens that for all microstates within the typical set,  $A(\xi \in \Omega_{typ}) = A^*$  (constant), then  $A$  is a state variable of the problem and  $\langle A \rangle = A^*$ .

Several generalizations of the typical set concept exist for example for a joint PDF in the joint typical set theorem<sup>37</sup>. The concept of typicality is extremely important also in thermodynamics. In the microcanonical ensemble where the PDF over the phase space is a uniform PDF, the entropy reduces to the integration of the density of accessible microstates.

However, as recent studies on stochastic thermodynamics have shown<sup>15,34</sup>, it is possible to extend many thermodynamic results also for small systems, i.e., systems where significant fluctuations of the state variables are not only possible, but also probable. Such type of systems are also those for which a practical implementation of a feedback controller is more feasible due to their limited number of degrees of freedom.

#### IV. LARGE DEVIATION THEORY

There is an elegant formalism within information theory connecting the statistical analysis of microscopic

fluctuations with the macroscopic behavior of the system described by thermodynamic potentials, this formalism is large deviation theory (LDT). This connection is a well established fact dating back to the '70, several authors<sup>38–43</sup> used LDT to derive many results of equilibrium thermodynamics. In particular it was possible to derive in a very elegant way the maximum entropy principle/ minimum free energy for systems in the thermodynamic limit in microcanonical or canonical ensemble.

The entire idea of LDT is to estimate the probability of fluctuations departing from the average in stochastic problems. This can be directly applied to evaluate the probability of observing a fluctuation of a state variable in a thermodynamic system.

Let us assume we have a quantity  $A$  with average value  $A^*$  and  $N$  degrees of freedom in the system. We define the contribution to  $A$  per degrees of freedom  $a = A/N$  and  $a^* = A^*/N$ . It is said that a stochastic problem follows a Large deviation (LD) law if the probability that  $a$  departs from  $a^*$  follows an exponential law:

$$q(a \neq a^*) \sim e^{-NK(a)}, \quad (12)$$

where  $K(a)$  is an exponent which is dependent on the thermodynamic quantity, while  $q(a)$  is defined as in eq. 11 from the PDF per microstate  $p(\xi)$ .

A powerful theorem to check if a problem satisfies a LD law is the Gärtner-Ellis theorem (GET)<sup>39</sup>. We first define the Scaled Cumulant Generating Function (SCGF) as

$$\lambda(\alpha) = \lim_{N \rightarrow \infty} \frac{1}{N} \ln [\langle e^{N\alpha a} \rangle], \quad (13)$$

with

$$\langle e^{N\alpha a} \rangle = \int p(\xi) e^{\alpha N A(\xi)} d\xi, \quad (14)$$

with  $\alpha$  a real number. The GET states that if the SCGF exists and is differentiable everywhere in  $\alpha$ , then the system fluctuations follow an exponential law.

We notice that the SCGF is very similar to a partition function, but scaled with respect to  $N$  and also where every exponential term is weighted by the probability per microstate,  $p(\xi)$ . If we make a simple variable change  $\alpha = -\beta$ , the new parameter  $\beta$  plays the same role as the inverse temperature,  $\beta = 1/(k_B T)$  and the SCGF can be rewritten as  $-\phi(\beta) = \lambda(\alpha)$ . The GET does not only provide a condition for existence, but also gives an operative way to evaluate the exponent  $K(a)$ . In fact it demonstrates that  $K(a)$  is related to the Legendre-Fenchel transform of the SCGF:

$$K(a) = -\min_{\beta \geq 0} [\beta a - \phi(\beta)]. \quad (15)$$

It is possible to demonstrate<sup>44</sup> that the previous equation can be rewritten in terms of the real partition function (without the PDF weighting the exponents as in the

SCGF), but adding a constant:

$$J(a) = -\min_{\beta \geq 0} [\beta a - \phi(\beta)] + \frac{1}{N} \ln \Lambda = K(a) + \frac{1}{N} \ln \Lambda. \quad (16)$$

The latter constant has the form of the entropy of a uniform PDF  $U(\xi) = 1/\Lambda$ , with  $\Lambda$  a particular volume within the phase space (see<sup>44</sup>), which depends on the prior PDF  $p(\xi)$ . In the case of a microcanonical ensemble it reduces to all the microstates with the same energy  $\bar{E}$ . The function  $\phi(\beta)$  is related to the free energy potential. Specifically, we have that the free energy per degrees of freedom ( $f = F/N$ ) is equal to:

$$f = \frac{\phi(\beta)}{\beta}, \quad (17)$$

thus the Legendre-Fenchel transform of the SCGF is linked to the entropy of the system per degree of freedom,  $s(a)$ , by

$$J(a) = -\frac{s(a)}{k_B} + \frac{1}{N} \ln \Lambda. \quad (18)$$

If we consider the constant term as the entropy associated to a uniform PDF we have that the exponent  $J(a)$  can be treated as follows:

$$\begin{aligned} J(a) &= \frac{-s(a)}{k_B} + \frac{1}{N} \ln \Lambda \\ &= \frac{1}{N} \left( \sum_{\xi} p(\xi) \ln p(\xi) + \sum_{\xi} p(\xi) \ln \Lambda \right) \\ &= \frac{1}{N} D(\Xi//U), \end{aligned} \quad (19)$$

where  $U = 1/\Lambda$  is the uniform PDF over the phase space volume  $\Lambda$ ,  $D$  is the KLd and  $Ns(a) = S_{tot} = -\sum p(\xi) \ln p(\xi)$  the total entropy. Thus the probability  $q(a \neq a^*)$  goes like the following:

$$q(a \neq a^*) \approx e^{-D(\Xi//U)}, \quad (20)$$

recovering the Chernoff bound for large fluctuations within the typical set formalism<sup>37</sup>. The LDT provides a clear understanding why entropy maximization is at the essence of equilibrium thermodynamics. Similar results are obtained for a canonical ensemble<sup>18</sup>. For a general discussion about the GET and the LDT applied to thermodynamics we refer to<sup>39</sup>.

Notably, LDT has been applied to non equilibrium systems<sup>38</sup> substituting the PDF for a microstate with the PDF of entire time dependent trajectories to take into account the time evolution of the system. Even more important, LDT can be used to derive fluctuation theorems. In fact if the exponent has the symmetry relation, e.g.  $K(-a) - K(a) = \gamma a$  ( $\gamma$  is a positive real number) for a certain thermodynamic quantity,  $A = Na$ , we immediately get the fluctuation theorem<sup>18,34,45-47</sup>:

$$\frac{q(a)}{q(-a)} \approx e^{N(K(-a) - K(a))} = e^{N\gamma a}. \quad (21)$$

A dual theorem of GET is the Varadhan theorem<sup>48</sup> which allows to invert the Legendre-Fenchel transform. This states that if

$$K(a) = -\min_{\beta \geq 0} [\beta a - \phi(\beta)], \quad (22)$$

is valid, then also the following relation holds:

$$\phi(\beta) = \min_a [\beta a + K(a)]. \quad (23)$$

The LD law in microcanonical or canonical ensemble explains, through the Varadhan theorem, why the minima of the free energy potentials are associated to the state variable values for the equilibrium state.

## V. RATE DISTORTION THEORY AND LARGE DEVIATION THEORY

In this paragraph we make explicit the connection between LDT and the information theory quantities presented in the previous sections and consider a system under feedback. In order to do so we need to introduce the rate distortion function (RDF), a central functional in rate distortion theory (RDT). RDT copes with a fundamental problem in communication, namely estimating the minimal information content that any message sent over a communication channel must contain such that the receiver can still have a good reconstruction of the original signal. The error source can be either due to distorting noise, a finite channel capacity or because of a lossy compression operated by the sender. In mathematical form, if we define a distance,  $d(\pi, \xi)$ , between the message sent,  $\xi$ , and the message received,  $\pi$  (eventually after decompression, noise deconvolution, etc...), the target of RDT is to find under which conditions the average distance can be kept lower than a given threshold,  $\Gamma$ <sup>37</sup>. Formally we ask,

$$\langle d(\xi, \pi) \rangle \leq \Gamma, \quad (24)$$

where the average is computed over the joint PDF  $p(\xi, \pi)$ . The distance,  $d$ , can be any functional with the properties of a distance (symmetry, positive definite,  $d = 0$  iff  $\xi = \pi$ , and must satisfy the Schwartz inequality). The most important theorem of RDT states that, given  $d$  and  $\Gamma$ , there exists a function,  $R(\Gamma)$  (the rate distortion function), representing the minimum information required in order to send messages with an average distortion not greater than  $\Gamma$ . The function  $R(\Gamma)$  has some remarkable properties: it is convex in the argument  $\Gamma$ , it converges to the entropy of the source (for a discrete PDF) for  $\Gamma = 0$ , while for  $\Gamma \geq \Gamma^*$  it is zero. In practice  $\Gamma^*$  is the limit distortion value after which the information content is completely lost and the receiver has equal chance by just guessing at random the most likely message, based on the joint PDF<sup>37</sup>. For an example of a RDF for a Gaussian PDF see Fig. 2.

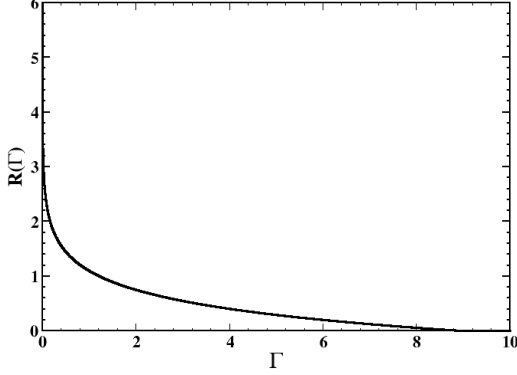


FIG. 2. Typical shape of a RDF. In this case it is plot the RDF of a Gaussian PDF with variance  $\sigma^2 = 9$ . For  $\Gamma$  larger than  $\sigma^2$  the RDF is 0.

The second important aspect of the rate distortion function is that it can be computed as a constrained minimization of a mutual information functional,

$$R(\Gamma) = \min_{p(\pi|\xi): \langle d(\xi, \pi) \rangle \leq \Gamma} I(\Xi, \Pi). \quad (25)$$

In the latter the minimization is with respect to the conditional PDF,  $p(\pi|\xi)$ , and the constrain is that the average distortion remains always smaller than  $\Gamma$ .

The connection between LDT and RDT is materialized by the distortion coding problem (DCP)<sup>49</sup>. DCP can be formulated in the following way: let assume we have two random vectors  $\xi$  and  $\pi$  with PDFs  $p(\xi)$  and  $q(\pi)$ , respectively, and we want to know what is the probability that picking up at random two vectors  $\xi$  and  $\pi$ , one from each distribution, the distance is such that  $d(\xi, \pi) \leq \Gamma$ . If we assume the only constrain that  $\xi$  should belong to the typical set of  $\Xi$ , the probability of such condition follows a LDT with an exponent equal to the rate distortion function:

$$p(\xi, \pi : \langle d(\xi, \pi) \rangle \leq \Gamma) \propto e^{-R(\Gamma)}. \quad (26)$$

The function  $R(\Gamma)$  monotonically decreases in the range  $0 < \Gamma < \Gamma^*$ , with limiting values  $R(0) = S/k_B$ ,  $R(\Gamma > \Gamma^*) = 0$ .

In the two limiting cases: when  $\Gamma = 0$ ,  $R(0) = S(\Xi)/k_B$ . In the case  $\Gamma \geq \Gamma^*$  the probability converges to 1 because  $R(\Gamma \geq \Gamma^*) = 0$ .

Since the rate distortion function appears in the form of a LD law, there is a second way to define the rate distortion function using the GET and the Legendre-Fenchel transform of the SCGF<sup>49,50</sup>:

$$R(\Gamma) = -\min_{\lambda \geq 0} \left[ \lambda \Gamma + \sum p(\xi) \ln(Z_\pi(\lambda)) \right], \quad (27)$$

with

$$Z_\pi(\lambda) = \sum_{\pi} q(\pi) e^{-\lambda d(\xi, \pi)}. \quad (28)$$

$Z_\pi$  has the form of a generalized partition function, linked to the distribution of  $\pi$ , where the distance  $d$  has a role similar to energy in the conventional partition function<sup>51</sup>.

An interesting aspect of this particular case of the LDT is that minimizing  $R(\Gamma)$  w.r.t.  $\lambda$ , we get an important relation,

$$\Gamma = \sum_{\xi, \pi} d(\xi, \pi) \frac{p(\xi)q(\pi)e^{-\lambda^* d(\xi, \pi)}}{Z_\pi(\lambda^*)} = \langle d(\xi, \pi) \rangle. \quad (29)$$

The average is made with respect to the joint PDF,

$$\tilde{p}(\xi, \pi) = \frac{p(\xi)q(\pi)e^{-\lambda^* d(\xi, \pi)}}{Z_\pi(\lambda^*)}, \quad (30)$$

which is the joint probability that fulfill the minimal rate function for a defined average error  $\Gamma$  and  $\lambda^*$  is the value minimizing the Legendre-Fenchel transform.

## VI. LARGE DEVIATION, RATE DISTORTION AND FEEDBACK CONTROL

We can now use the results of RDT applied to LDT to obtain our most important result concerning the thermodynamics of systems with feedback control. Let assume that we have a feedback controller that performs measurements of the state of a system and then afterwards manipulates it. If we assume that the measurement,  $\pi_k$ , has some correlation with the state  $\xi$ , then the entropy of the system after measurement will be given by the PDF,  $p(\xi|\pi = \pi_k)$ , conditioned by the outcome  $\pi_k$ . This is equal to:

$$S(\Xi|\pi = \pi_k) = -k_B \sum p(\xi|\pi = \pi_k) \ln p(\xi|\pi = \pi_k). \quad (31)$$

A typical set volume is always associated to this entropy,

$$\Omega_{typ}(\xi|\pi = \pi_k) = e^{\frac{S(\Xi|\pi = \pi_k)}{k_B}}, \quad (32)$$

hence the average volume of the typical set after a measurement is thus:

$$\langle \Omega_{typ}(\xi|\pi = \pi_k) \rangle = \left\langle e^{\frac{S(\Xi|\pi = \pi_k)}{k_B}} \right\rangle, \quad (33)$$

where the average is made w.r.t.  $q(\pi)$ .

The lower boundary of this formula can be obtained using the Jensen inequality, thanks to the convexity of the exponential function,  $\langle \exp(f) \rangle \geq \exp(\langle f \rangle)$ , to obtain:

$$\langle \Omega(\xi|\pi = \pi_k)_{typ} \rangle \geq e^{\frac{\langle S(\Xi|\pi = \pi_k) \rangle}{k_B}} = e^{\frac{S(\Xi|\Pi)}{k_B}}, \quad (34)$$

exploiting the fact that  $\langle S(\Xi|\pi = \pi_k) \rangle = S(\Xi|\Pi)$ .

Thus, we find that the feedback can -at best- constrain the volume of phase-space by the conditional entropy. Splitting the conditional entropy in the usual,  $S(\Xi|\Pi) = S(\Xi) - k_B I(\Xi, \Pi)$ , we obtain that the effect

of the feedback is to compress the volume of the typical set with an exponent at best equal to the mutual information:

$$\Omega_{typ}(\xi|\pi) \geq \Omega_{typ}(\xi)e^{-I(\Xi, \Pi)}. \quad (35)$$

The fact that the mutual information is always non negative ensures that the effect of the feedback is always to compress the original typical set.

If we want to connect the mutual information to the effective measurement operated by the feedback apparatus (FA), then we can define a distance functional,  $d$ , and an average distortion  $\Gamma$  between state and estimation and use the rate distortion function to find the minimal mutual information required to have a certain average distortion  $\Gamma$ . The final result is the lower boundary,

$$\langle \Omega_{typ}(\xi|\pi = \pi_k) \rangle \geq \Omega_{typ}(\xi)e^{-R(\Gamma)}. \quad (36)$$

The appealing of equation (36) is that it gives a direct connection between the effect of feedback information and the effective operative measurement performed by the FA.

A final note regarding the effect of the FA. The reduction in the volume of the typical set remains only "virtual" until the FA does not operate and manipulate the system. We can thus imagine this shrinking of the typical set linked to the  $R(\Gamma)$  as the best average reduction in uncertainty about the state of the system from the FA side once the measurements have been made. It is anyway clear that the entropy and the physics of the system are left unchanged until the FA does not directly operate.

## VII. TOWARDS A THERMODYNAMIC POTENTIAL INCLUDING INFORMATION

A thermodynamic potential is nothing else than a quantity that when minimized/maximized allows to find the equilibrium or stationary states of a certain thermodynamic system, subject to given constraints. We have already discussed in section 4 and 5 how the LDT provides an elegant way to derive the maximal entropy and minimal energy principles for the microcanonical and canonical ensembles. In particular, we saw how the free energy potential comes as part of the Legendre-Fenchel transform of the large deviation exponent thanks to the Gartner-Ellis theorem. In this section we put together this result and the results of section 6 in order to construct an explicit thermodynamic potential for a system under feedback control.

First of all we define a simplified feedback apparatus (SFA). We assume a system in thermal equilibrium with an external bath at temperature  $T_0$  and we assume that the system can be coupled to the bath or disconnected from that and coupled to the SFA. The SFA is a system that can make a measurement and manipulate the

system accordingly in a time  $\tau$  much smaller than any relaxation time of the system. Moreover, we assume a certain level of ideality in the feedback, i) during the measurement no perturbation of the system occurs, ii) the feedback uses the entire information during the manipulation achieving the maximal efficiency. Once the system has been manipulated by the feedback it is allowed to relax and finally it is connected again to the external bath until thermalize with it. This simple cycle assures that every time the feedback performs its new measurement/manipulation the system is at thermal equilibrium with temperature  $T_0$ .

With these approximations we can decouple the initial state before the feedback operation from the feedback action. A completely different and more complex scenario occurs in the case when the system is not allowed to thermalize or the feedback acts continuously in such a way that the system state depends on previous feedback history.

We finally assume that the feedback can obtain an estimate,  $\pi$ , of the real state,  $\xi$ , of the system such that  $\langle d(\xi, \pi) \rangle \leq \Gamma$  for a well-defined distance functional.

We have shown in section 6 that the feedback action entails to an entropy reduction after measurement by a factor  $I(\xi, \pi) = R(\Gamma)$ . Thus we can finally define the maximum increase in free energy due to the presence of the feedback:

$$\begin{aligned} \Delta F &= F - F_{eq} \\ &= U - T_0(S - k_B R(\Gamma)) - U + T_0 S = k_B T_0 R(\Gamma) = \langle W \rangle. \end{aligned} \quad (37)$$

with  $U$  the internal energy. It recovers the result found by Sagawa with eq. 1, as the maximum work for a cyclic transformation, but now with a direct connection between the type and accuracy of the measurement,  $d(\xi, \pi)$  and  $\Gamma$ , and the average extracted work  $\langle W \rangle$ .

The  $R(\Gamma)$  exponent of equation (36), following a LD law, can be expanded in the Legendre-Fenchel transform as in equation (15), leading to a joint probability between the microstate and the guess equal to equation (29). Using the Varadhan theorem we can invert the Legendre-Fenchel transform for the RDF and obtain something equivalent to a thermodynamic potential for the information:

$$\langle \ln(Z_\pi(\lambda)) \rangle = \sum p(\xi) \ln(Z_\pi(\lambda)) = \min_{\Gamma} [\lambda \Gamma + R(\Gamma)]. \quad (38)$$

The structure of this equation is similar to the one of the free energy in standard thermodynamics with the relation between free energy and partition function, where now the Chernoff coefficient  $\lambda$  plays the role of inverse temperature:

$$\phi_I(\lambda) \equiv -\frac{1}{\lambda} \langle \ln(Z_\pi(\lambda)) \rangle. \quad (39)$$

What is the meaning of this functional? We can understand its role by calculating the distance, using the

KLd, between a generic joint probability  $p(\xi, \pi)$  and the optimal  $\tilde{p}(\xi, \pi)$  as defined in eq. (29):

$$\begin{aligned} D(P//\tilde{P}) &= \sum_{\xi\pi} p(\xi, \pi) \ln \frac{p(\xi, \pi)}{\tilde{p}(\xi, \pi)} \\ &= \sum_{\xi\pi} p(\xi, \pi) \log \frac{p(\xi, \pi)}{p(\xi)q(\pi)} + \lambda^* \langle d(\xi, \pi) \rangle + \langle \ln Z_\pi(\lambda^*) \rangle \\ &= I(\Xi, \Pi) + \lambda^* \langle d(\xi, \pi) \rangle - \lambda^* \phi_I(\lambda^*). \end{aligned} \quad (40)$$

The latter equation can be rewritten as:

$$\lambda^* \phi_I(\lambda^*) = \lambda^* \langle d(\xi, \pi) \rangle + I(\Xi, \Pi) - D(P//\tilde{P}). \quad (41)$$

Considering that the three terms in the rhs are all positive it is easy to demonstrate that the maximum of the potential  $\phi_I(\lambda^*)$  is obtained for  $D(P//\tilde{P}) = 0$ , that is when the joint probability is the ideal one for which the RDF is achieved.

### VIII. A SIMPLE PHYSICAL MODEL

We apply the formalism to a simple model: a set of single particles in a box in gas phase. We assume that every particle is under a feedback and that it is in thermal equilibrium with a bath at temperature  $T$ . The Hamiltonian of the system is given by  $\mathcal{E} = Ap^2$ , where  $A = 1/2m$ , being  $m$  the particle mass. The particle state,  $\xi = (r, p)$ , is characterized by an homogeneous spatial distribution for the position  $r$  within the box, and a normal distribution of the momentum,  $p$ . Therefore, neglecting position, we identify the particle state just with the momentum,  $x = p$ . The equilibrium distribution is a normal distribution with 0 mean value and a variance  $\sigma_\xi^2 = k_B T / 2A$ . The model includes a feedback that can probe the momentum of the particle. The feedback uses a distance  $d(\xi, \pi) = (\xi - \pi)^2$ . This measurement is affected by error, so the measurement  $\pi$  is a random variable, statistically correlated with the dynamical state of the particle. The model assumes that between every probe and manipulation the system is allowed to relax to thermal equilibrium, that means that at every measurement the system is found in the same equilibrium state. Finally, we also assume that  $\pi$  is distributed like a normal random variable. The situation can be formalized like in<sup>37</sup> by assuming that the feedback and the source are connected by a channel with Gaussian noise (see Fig. 3).

The Gaussian noise has distribution:

$$p(z) = N(0, \Gamma) = \frac{1}{\sqrt{2\pi\Gamma}} e^{-\frac{z^2}{2\Gamma}}. \quad (42)$$

For simplicity we assume that the source has mean value equal to zero ( $\mu_\xi = 0$ ),  $p(\xi) = N(0, \sigma_\xi^2)$ . This assumption is totally immaterial for the generality of the discussion.

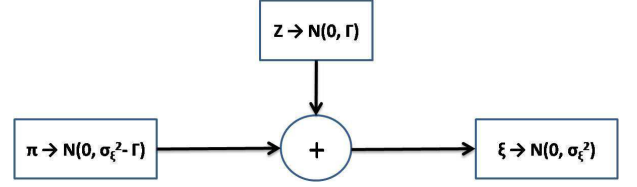


FIG. 3. (color online) Scheme of a channel between two joint gaussian random variables.

Following<sup>37</sup> the feedback has a PDF equals to:

$$p(\pi) = N(0, \sigma_\xi^2 - \Gamma) = \frac{1}{\sqrt{2\pi(\sigma_\xi^2 - \Gamma)}} e^{-\frac{\pi^2}{2(\sigma_\xi^2 - \Gamma)}}. \quad (43)$$

In the latter we have assumed that  $0 \leq \Gamma \leq \sigma_\xi^2$ . With the previous defined distance functional the average distance  $\Gamma$  is equal to the mean square error  $\langle (\xi - \pi)^2 \rangle$ .

For this simple model the form of the rate distortion function is well known:

$$R(\Gamma) = \frac{1}{2} \ln \left( \frac{\sigma_\xi^2}{\Gamma} \right), \quad (44)$$

with  $R(\Gamma) = 0$  for  $\Gamma \geq \sigma_\xi^2$ .

It is a simple matter of calculation to demonstrate that indeed eq. 27, using  $p(\xi)$ ,  $p(\pi)$  and  $d(\xi, \pi) = (\xi - \pi)^2$ , gives back the correct rate distortion function. The relation between  $\lambda$  and  $\Gamma$  is very simple:

$$\lambda = \frac{1}{2\Gamma}. \quad (45)$$

We can insert the PDF and the distance functional in eq. 30 in order to get the optimal joint distribution for such marginal PDFs. The solution is a joint Gaussian distribution:

$$\begin{aligned} \tilde{p}(\xi, \pi) = & \frac{1}{2\pi\sigma_\xi\sigma_\pi\sqrt{1-\rho^2}} \exp \left( -\frac{1}{2(1-\rho^2)} \left[ \frac{\xi^2}{\sigma_\xi^2} - \frac{2\rho\xi\pi}{\sigma_\xi\sigma_\pi} + \frac{\pi^2}{\sigma_\pi^2} \right] \right) \end{aligned} \quad (46)$$

with correlation coefficient equal to  $\rho = \sqrt{1 - \Gamma/\sigma_\xi^2}$ . Using the result in eq. 45 we obtain the relation between  $\lambda$  and the correlation coefficient:

$$\lambda = \frac{1}{2\sigma_\xi^2(1-\rho^2)}. \quad (47)$$

As expected if  $\Gamma$  is related to the distance between state and estimation by the feedback,  $\lambda$  has a simple interpretation in terms of correlation between the two random variables. In case of random vectors (with  $N$  elements) the total average distance  $\Gamma = \sum N\Gamma_i$ , with  $\Gamma_i$  the distance associated to the  $i^{th}$  component. Clearly the value

$\Gamma$  has an extensive characteristic being function of the number of degree of freedom ( $N$ ) of the system. On the contrary  $\lambda$  is the intensive counterpart being connected to the correlation between state and estimation.

In this particular case we can also calculate the optimal work that such a feedback apparatus can recover from the system as a function of the average distance  $\Gamma$ , using eq.37:

$$\langle W \rangle = Tk_B R(\Gamma) = -\frac{Tk_B}{2} \ln(1 - \rho^2) = Tk_B I(\Xi, \Pi). \quad (48)$$

The latter shows that the maximum work extracted is proportional as should be to the mutual information between two Gaussian random variables.

## IX. CONCLUSIONS

In this work we have analyzed a system under feedback using typicality and large deviation theory. In particular this connection assures to make a simple and natural link between the way the measurement is performed and its

accuracy and the maximum work that can be extracted in terms of entropy reduction. Clearly in this discussion we have not considered the effect of the manipulation from the feedback and its cost in terms of overall efficiency, thus all our results must be considered as boundary limits.

The main result of the paper is the possibility, using rate distortion theory, of developing the equivalent of thermodynamic potentials even in the case of systems under an external feedback.

This perspective to the problem is particular appealing not only because it establishes a nice relation between the measurement and the respective entropy and information functional, but also because it can be generalized for a large class of problems where large deviation applies.

## X. ACKNOWLEDGEMENTS

We acknowledge Prof. Merhav for the useful comments.

- 
- <sup>1</sup> H. S. Leff, A. F. Rex, "Maxwell Demon: Entropy, Information, Computing, Princeton University Press", Princeton, NJ, (1990).
  - <sup>2</sup> H. S. Leff, A. F. Rex, "Maxwell Demon 2: Entropy, Classical and Quantum Information, Computing", Institute of Physics, Bristol, (2003).
  - <sup>3</sup> L. Szilard, On the Decrease of Entropy in a Thermodynamic System by the Intervention of Intelligent Beings, *Z. Phys.* 53:840 (1929): English translation reprinted *Behavioral Science* 9:301 (1964).
  - <sup>4</sup> H. Touchette and S. Lloyd, *Phys. Rev. E*, 84, 1156 (2000).
  - <sup>5</sup> A. E. Allahverdyan, D. Janzing, and G. Mahler, *J. Stat. Mech.*, P09011 (2009).
  - <sup>6</sup> J. Horowitz, T. Sagawa, and J. M. R. Parrondo, *Phys. Rev. Lett.*, 111, 010602 (2013).
  - <sup>7</sup> J. M. Horowitz and M. Esposito, *Phys. Rev. X*, 4, 031015 (2014).
  - <sup>8</sup> J. M. Horowitz and H. Sandberg, *New J. Phys.*, 16, 125007 (2014).
  - <sup>9</sup> N. Shiraishi, S. Ito, K. Kawaguchi, and T. Sagawa, *New J. Phys.*, 17, 045012 (2015).
  - <sup>10</sup> D. Hartich, A. C. Barato, and U. Seifert, *J. Stat. Mech.*, P02016 (2014).
  - <sup>11</sup> S. Toyabe, T. Sagawa, M. Ueda, E. Muneyuki and M. Sano, *Nature Physics*, 6, 988 (2010).
  - <sup>12</sup> T. Sagawa and M. Ueda, *Phys. Rev. Lett.*, 100, 080403 (2008).
  - <sup>13</sup> T. Sagawa and M. Ueda, *Phys. Rev. Lett.*, vol. 104, 090602 (2010).
  - <sup>14</sup> T. Sagawa and M. Ueda, Chap. 6, "Nonequilibrium Statistical Physics of Small systems", (Wiley-VCH), (2013).
  - <sup>15</sup> S. Hilbert, P. Haenggi and J. Dunkel, *Phys. Rev. E*, vol. 90, 062116 (2014).
  - <sup>16</sup> J. M. R. Parrondo, J. M. Horowitz and T. Sagawa, *Nature Phys.*, vol. 11, 131 (2015).
  - <sup>17</sup> A. Gagliardi and A. Pecchia, arXiv:1503.02824v1 (2015).
  - <sup>18</sup> H. Touchette, *Phys. Rep.*, vol. 478, 1-69 (2009).
  - <sup>19</sup> N. Merhav, *IEEE Trans. Inform. Theory*, vol. 54, no. 8, pp. 3710-3721 (2008).
  - <sup>20</sup> A. Kis Andras and A. Zettl, *Philos. Trans. R. Soc. A*, 366, 1591 (2008).
  - <sup>21</sup> J. R. Gomez-Solano, L. Bellon, A. Petrosyan and S. Ciliberto, *Europhys. Lett.*, 89, 60003 (2010).
  - <sup>22</sup> L. Bellon, L. Buisson, S. Ciliberto and F. Vittoz, *Rev. Sci. Instrum.*, 73, 3286 (2002).
  - <sup>23</sup> A. Berut, A. Arakelyan, A. Petrosyan, S. Ciliberto, R. Dillenschneider and E. Lutz, *Nature*, 187, 483 (2012).
  - <sup>24</sup> J. V. Koski, V. F. Maisi, T. Sagawa and J. P. Pekola, *Phys. Rev. Lett.*, 113, 030601 (2014).
  - <sup>25</sup> A. C. Barato, D. Hartich and U. Seifert, *Phys. Rev. E*, 87, 042104 (2013).
  - <sup>26</sup> S. Ito, T. Sagawa, *Nat. Comm.*, 6, 7498 (2015).
  - <sup>27</sup> A. H. Lang, C. K. Fisher, T. Mora and P. Mehta, *Phys. Rev. Lett.*, 113, 148103 (2014).
  - <sup>28</sup> P. Sartori, L. Granger, C. F. Lee and J. M. Horowitz, *PLoS Comput. Biol.*, 10, 1003974 (2014).
  - <sup>29</sup> R. G. Endres and N. S. Wingreen, *Phys. Rev. Lett.*, 103, 158101 (2009).
  - <sup>30</sup> H. Qian and T. C. Reluga, *Phys. Rev. Lett.*, 94, 028101 (2005).
  - <sup>31</sup> P. Mehta and D. Schwab, *Proc. Natl Acad. Sci. USA*, 109, 17978 (2012).
  - <sup>32</sup> Y. Tu, *Proc. Natl Acad. Sci. USA*, 105, 11737 (2008).
  - <sup>33</sup> N. Merhav, *J. Stat. Phys.*, P01029, doi: 10.1088/1742-5468/2011/01/P01029 (2011).
  - <sup>34</sup> U. Seifert, *Rep. Prog. Phys.*, 75, 126001 (2012).



- <sup>35</sup> C. Jarzynski, Non-equilibrium equality for free energy differences, *Phys. Rev. Lett.*, vol. 78, 2690 (1997).
- <sup>36</sup> F. Rezakhanlou and C. Villani, "Entropy Methods for the Boltzmann Equation", Springer-Verlag Berlin Heidelberg (2008).
- <sup>37</sup> T. M. Cover and J. A. Thomas, "Elements of information theory", Wiley (2006).
- <sup>38</sup> H. Touchette, R. J. Harris, Chap. 11, "Nonequilibrium Statistical Physics of Small systems", (Wiley-VCH), (2013).
- <sup>39</sup> R. S. Ellis, "Entropy, Large Deviations and Statistical Mechanics", Springer, New York (1985).
- <sup>40</sup> R. S. Ellis, *Physica D*, 133, 106-136 (1999).
- <sup>41</sup> Y. Oono, *Progr. Theoret. Phys. Suppl.* 99, 165-205 (1989).
- <sup>42</sup> R. S. Ellis, *Scand. Actuar. J.*, 1, 97-142 (1995).
- <sup>43</sup> O. E. Lanford, "Entropy and equilibrium states in classical statistical mechanics", in: A. Lenard (Ed.), *Statistical Mechanics and Mathematical Problems*, vol. 20, Springer, Berlin, pp. 1-113 (1973).
- <sup>44</sup> N. Merhav, *IEEE IST (2008) Toronto*, 499 (2008).
- <sup>45</sup> D. J. Evans and D. J. Searles, *Phys. Rev. E*, 50, 1645 (1994).
- <sup>46</sup> G. Gallavotti and E. G. D. Cohen, *Phys. Rev. Lett.*, 74, 2694 (1995).
- <sup>47</sup> J. L. Lebowitz and H. Spohn, *J. Stat. Phys.*, 95, 333 (1999).
- <sup>48</sup> S. R. S. Varadhan, *Comm. Pure App. Math.*, 19, 261 (1966).
- <sup>49</sup> N. Merhav, *Statistical Physics and Information theory, Foundation and Trends in Communication and Information Theory*, vol. 6, 1-212 (2009).
- <sup>50</sup> T. Berger, "Rate distortion theory: a mathematical basis for data compression", PrenticeHall, Inc., Engelwood Cliffs, NJ, (1971).
- <sup>51</sup> R. M. Gray, "Source Coding Theory", Kluwer Academic Publishers, (1990).